# High-dimensional data: introduction

Wessel van Wieringen
w.n.van.wieringen@vu.nl

Department of Epidemiology and Biostatistics, VUmc
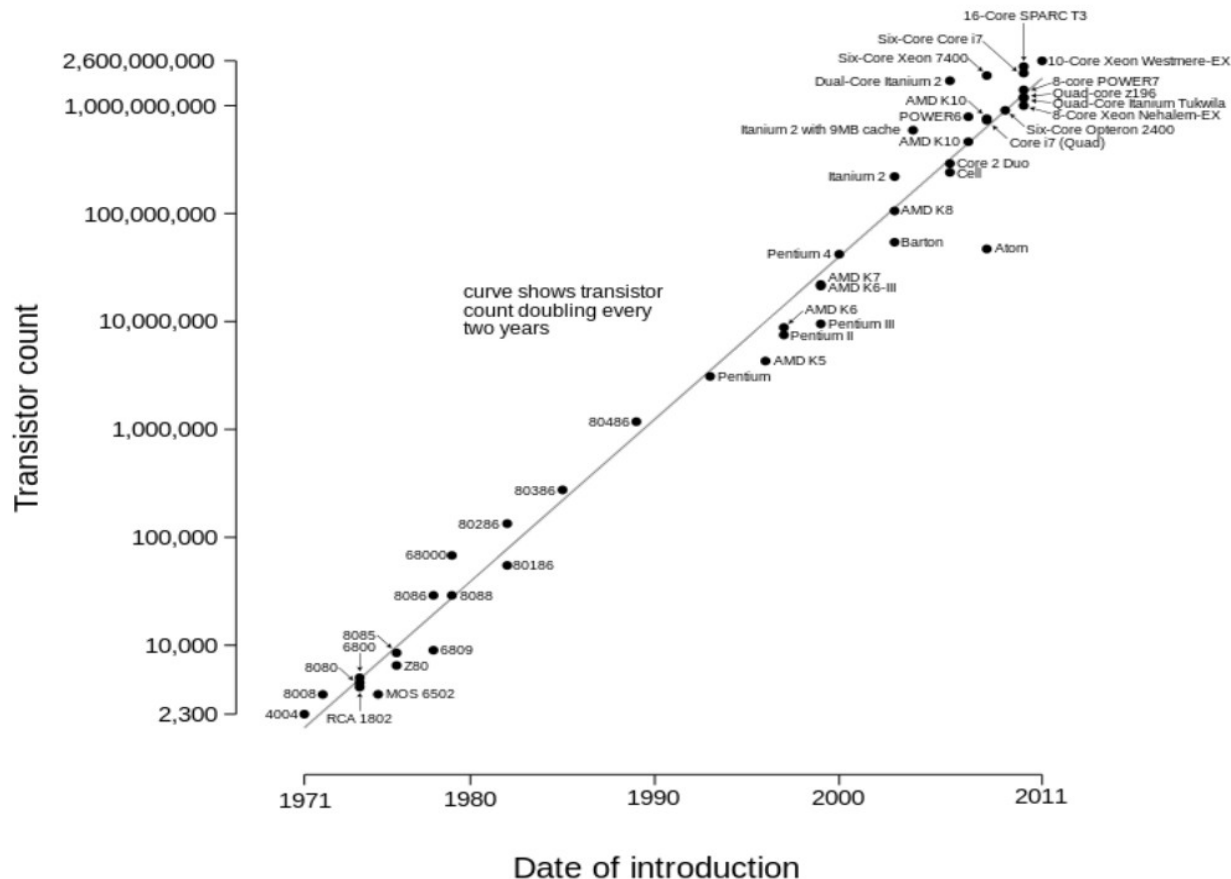& Department of Mathematics, VU University
Amsterdam, The Netherlands

vrije Universiteit

VU medisch centrum
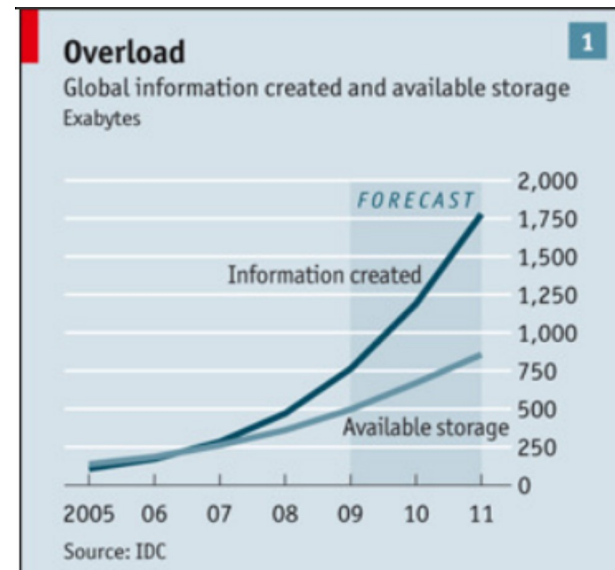
# How did we end up here?

*Moore's law*

The number of transistors in a dense integrated circuit doubles approximately every two years.



source: wikipedia

# How did we end up here?

*Data deluge*

"... the quantity of information in the world is soaring. According to one estimate, mankind created 150 exabytes (billion gigabytes) of data in 2005. This year, it will create 1,200 exabytes. Merely keeping up with this flood, and storing the bits that might be useful, is difficult enough. Analysing it, to spot patterns and extract useful information, is harder still."
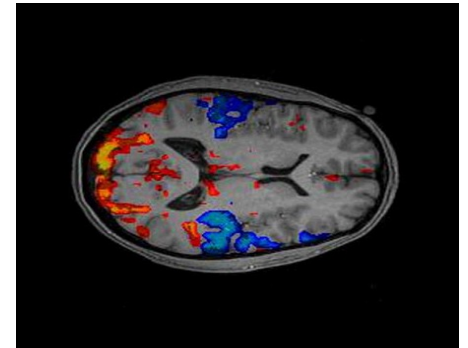


**Overload**
Global information created and available storage
Exabytes

FORECAST

Information created

Available storage

2,000
1,750
1,500
1,250
1,000
750
500
250
0

2005 06 07 08 09 10 11

Source: IDC
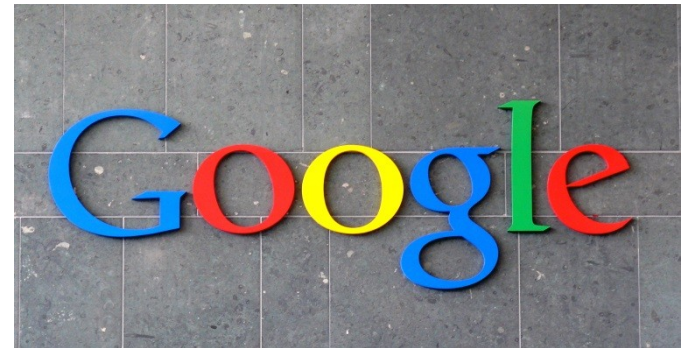
# How did we end up here?

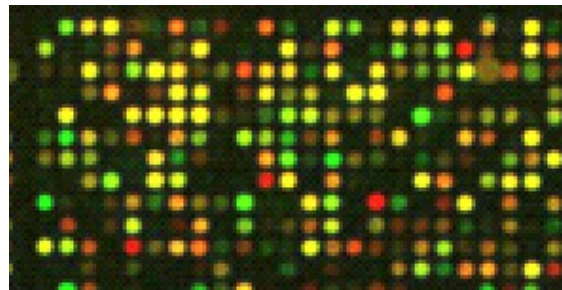*Examples*

→ Brain image data (fMRI / EEG)



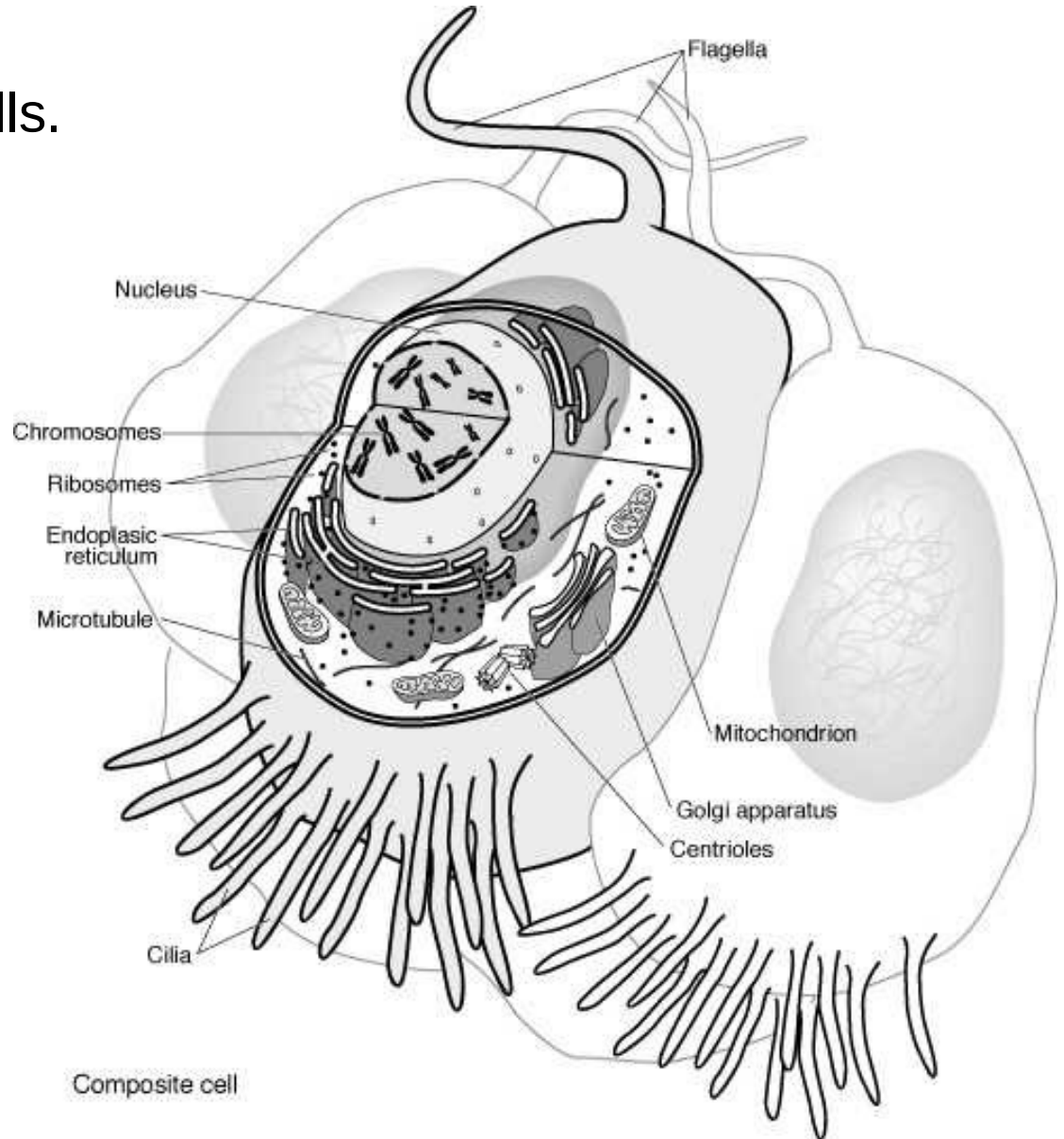→ Movie database



→ Google search data



→ Microarrays

# A minimum of biology

Organisms are made of cells.

A *cell* is the smallest possible independent living unit. The cell contains a complete copy of the organisms genome.
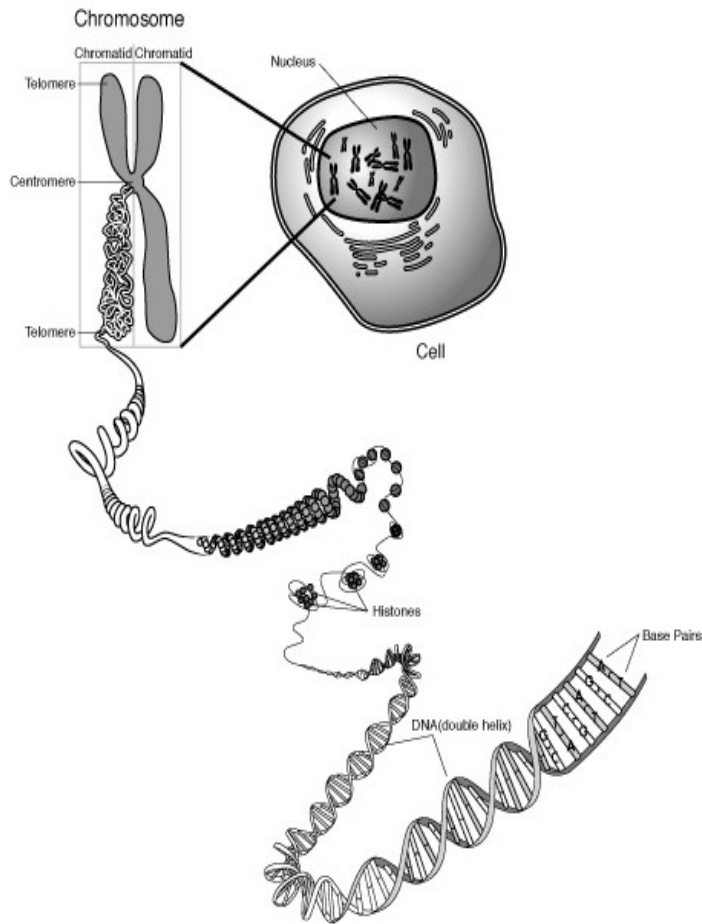
The *genome* is the total genetic constitution of an organism, the full haploid set of chromosomes with all its genes.

Flagella

Nucleus

Chromosomes

Ribosomes

Endoplasic reticulum

Microtubule

Mitochondrion

Golgi apparatus
Centrioles

Cilia

Composite cell

# A minimum of biology

A *chromosome* is one of a set of threadlike molecular structures composed of compressed *DNA*, that carry the genes which determine an individual's heriditary traits.
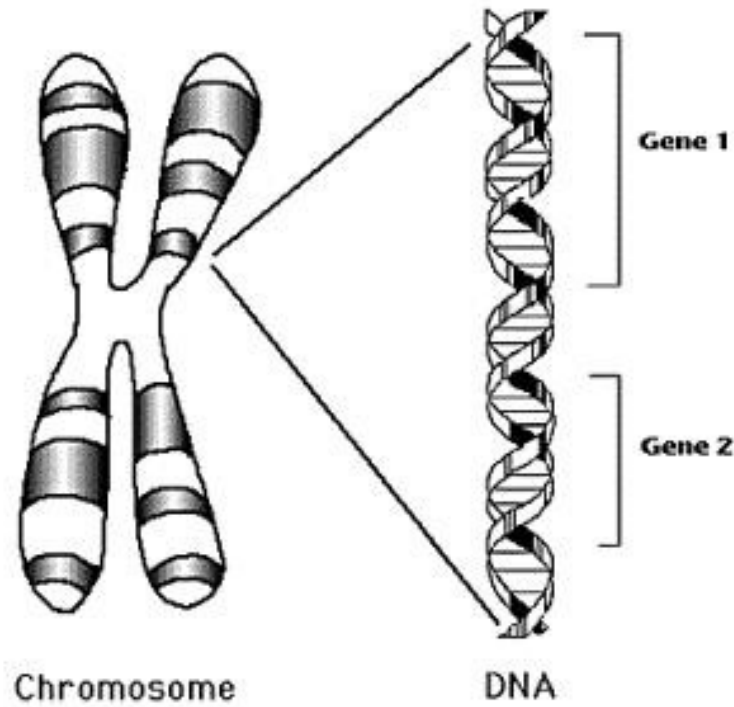


Conceptually, *DNA* is an information carrier, information necessary for the functioning of cells and encoded in molecular units called genes.

On the molecular level *DNA* is a double-stranded polymer composed of four basic molecular units called nucleotides.

# A minimum of biology

A *gene* is the basic physical unit of heredity: a linear sequence of nucleotides, as a segment of DNA located on a chromosome, that provides the coded instruction for one polypeptide chain.
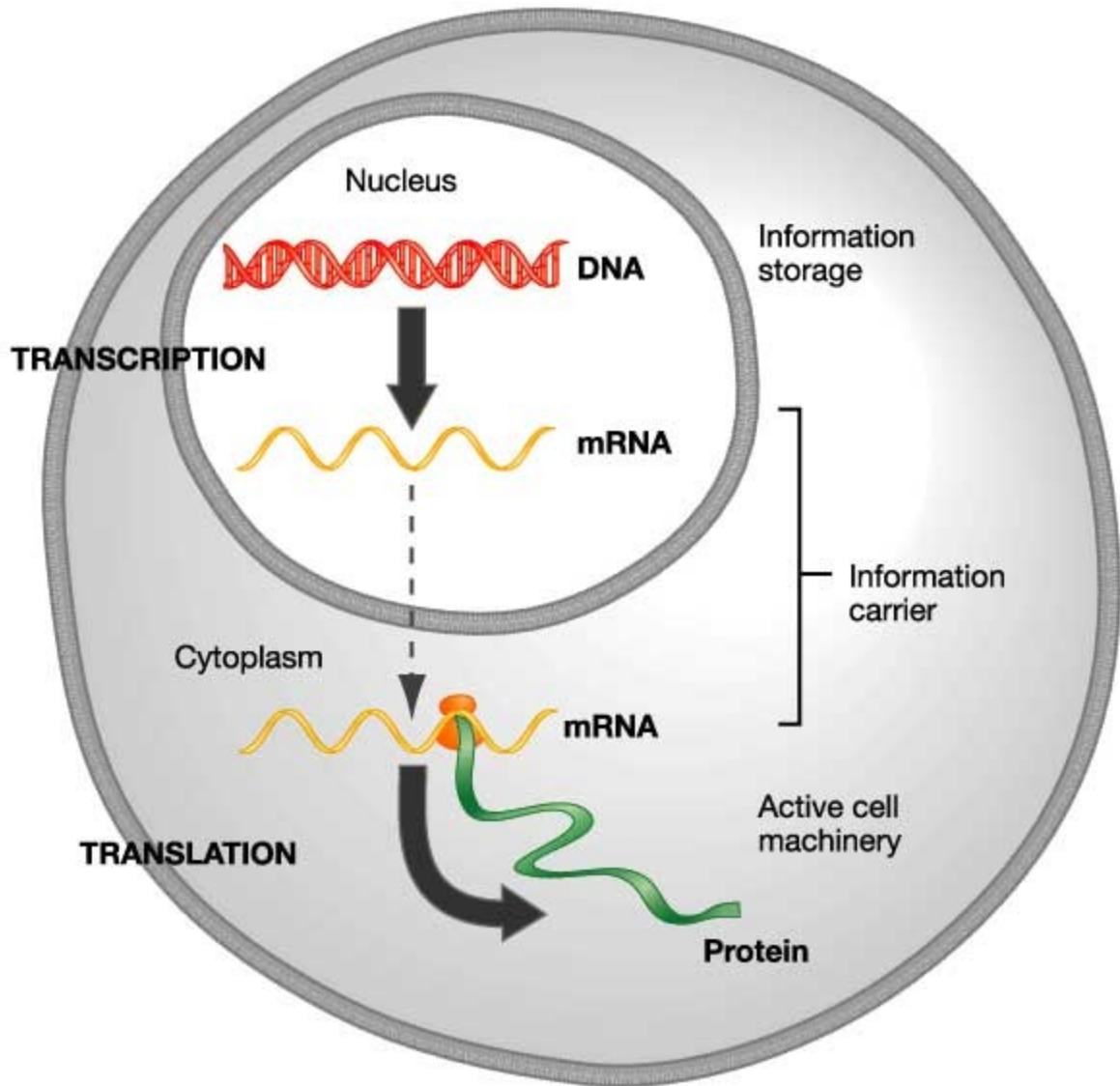


Chromosome          DNA

Gene 1

Gene 2

# A minimum of biology

*Central Dogma of Molecular Biology* describes the information transfer process that leads from the information encoded in DNA to the proteins in the cell.

Three steps are discerned:
1) Replication
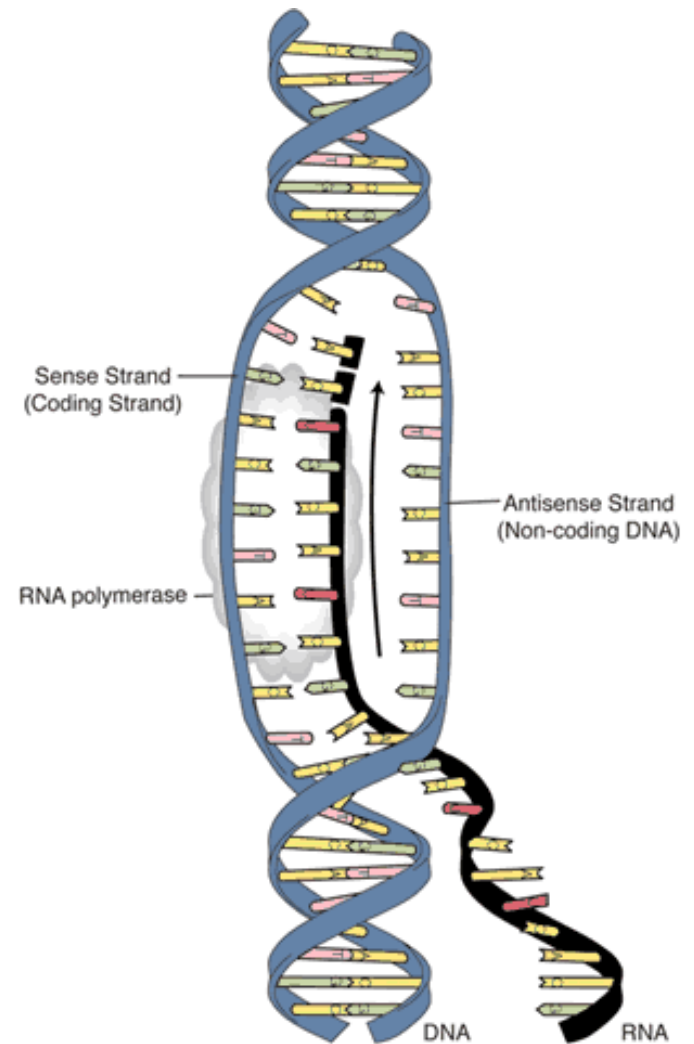2) Transcription
3) Translation

# A minimum of biology

*Transcription* is the synthesis of *mRNA* (nucleic acid like DNA, but single-stranded) from the DNA in the nucleus.

The mRNA is transported to the cytoplasm and used to synthesize protein.

*Jargon*
A gene is said to be *expressed* if the product it encodes for has been formed.



Sense Strand
(Coding Strand)

Antisense Strand
(Non-coding DNA)

RNA polymerase

DNA          RNA

# A minimum of biology

Molecular biology aims to understand the molecular processes that occur in the cell. That is, which molecules present in the cell interact, and how is this coordinated?

For many cellular process, it is unknown which genes play what role.

*Solution*
Simply measure (the expression of) all genes ...
… and later sort out which are relevant.

# Microarrays

*Microarray*

- Conceptually: a measurement device.

Gene expression arrays measure the expression of genes (which genes are expressed and to what extent).

In fact, it measures mRNA which is related – through the transcription process – to the expression of genes.

Other types of microarrays measure:
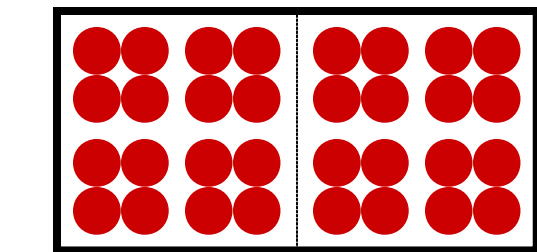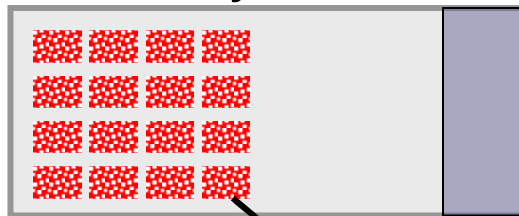
- SNPs

- DNA copy number
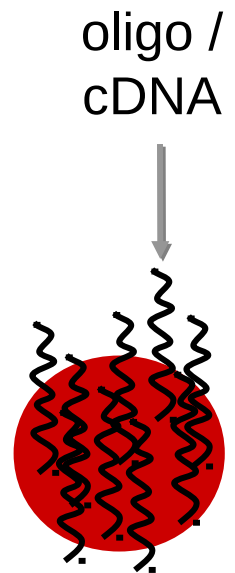
- methylation

- ...

# Microarrays

## *Microarray*

- Physically: a glass slide


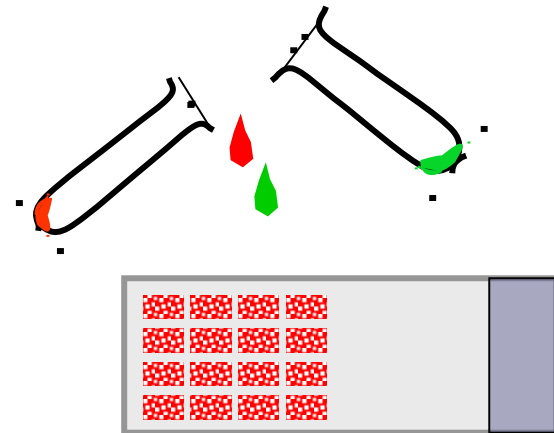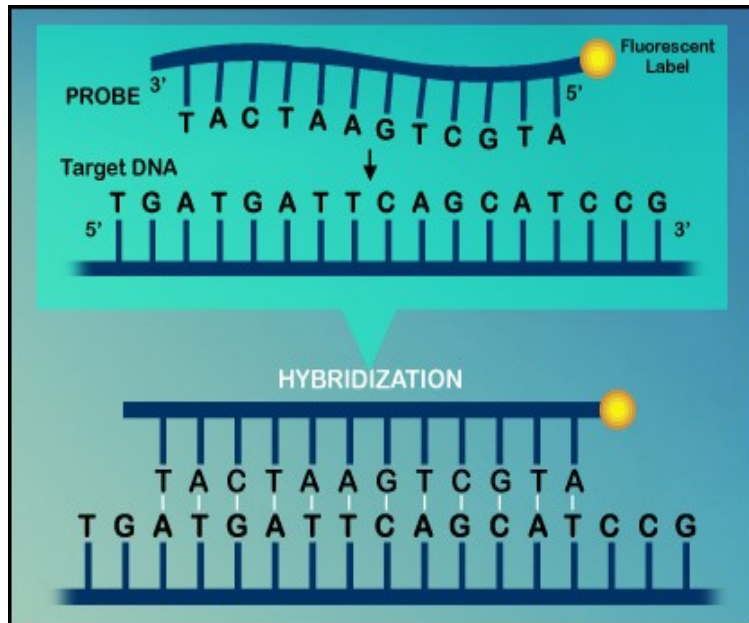
Microarray



Collection of features / probes

oligo / cDNA

Feature / probe

# Microarrays

*Hybridization*
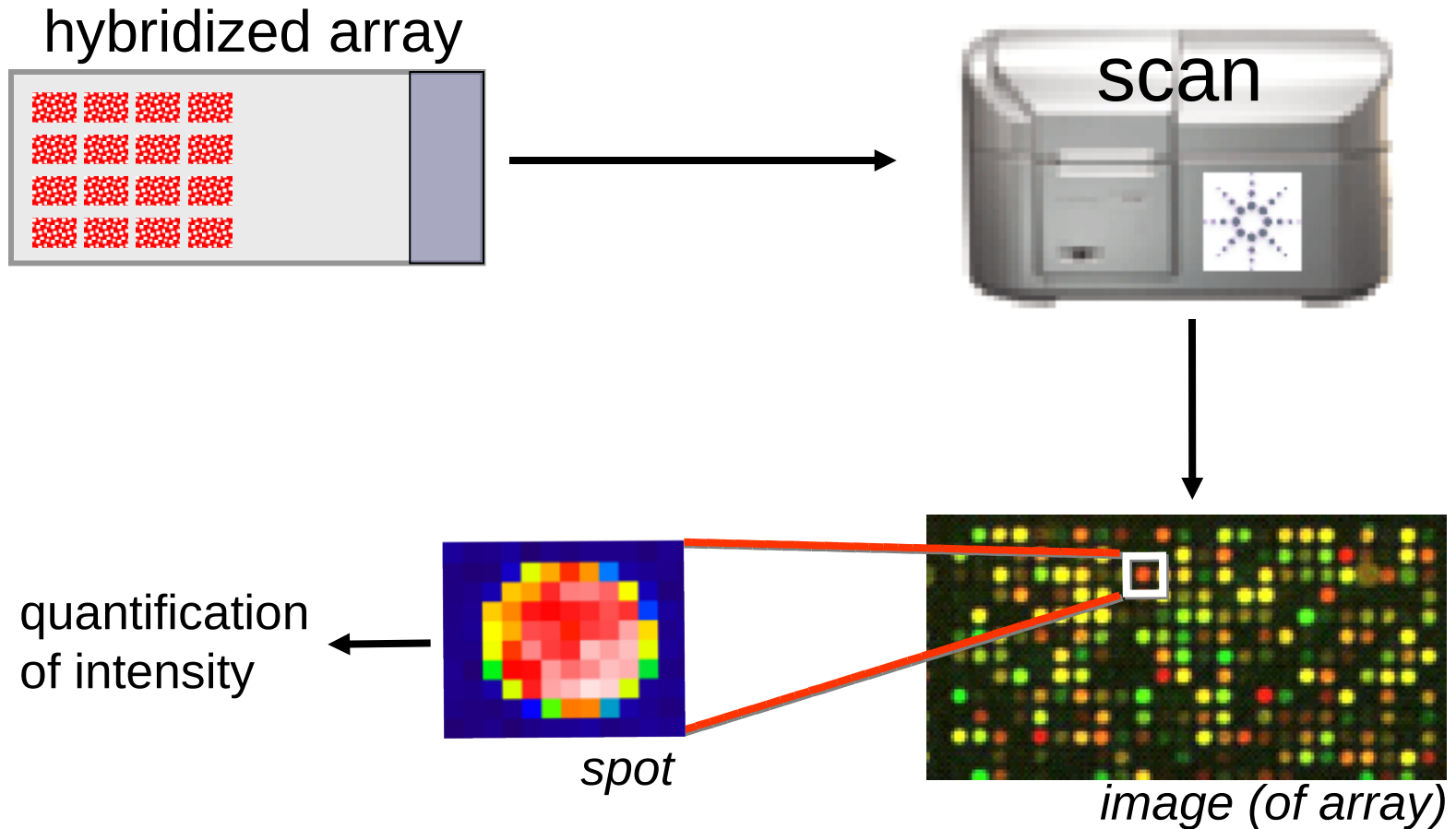
The preferential binding of gene sequences to complementary sequences.

# Microarrays

An image of the microarray is generated

hybridized array

scan

quantification
of intensity

*spot*

*image (of array)*

# Microarrays

Output per hybridization:
- File of 44Mb
- ~44000 rows
- ~100 columns
- Annotation information
- Quality metrics
- Biological signal (in various forms)
- Background signal (in various forms)

# Quality control

Plot the raw image of the array

```
> image(…)
```



Nothing special



…

# Quality control

Plot image of fore- and background signal

```
> image(…)
```



Foreground signal



Background signal

# Quality control

Generate boxplots of fore- and background signal

```
> boxplot(…)
```



Foreground

Background

# Preprocessing

Before the statistical analysis of interest, the gene expression measurements (intensities) undergo several preprocessing steps.

```
┌─────────────────────────────────┐
│      Background correction       │
└─────────────────────────────────┘
                 │
                 ▼
┌─────────────────────────────────┐
│    Within-array normalization    │
└─────────────────────────────────┘
                 │
                 ▼
┌─────────────────────────────────┐
│   Between-array normalization    │
└─────────────────────────────────┘
                 │
                 ▼
┌─────────────────────────────────┐
│       Expression measure         │
└─────────────────────────────────┘
```

# Background correction

*Background intensity*
The background of the microarray may have a non-zero intensity.

Hence, a feature's intensity may include a contribution not specifically due to the hybrization of the target to the probe.



intensity

measured signal

≈ noise

spot          background

spot



spot surrounding



nonzero intensity: not necessary representative of background noise inside feature area.

signal



contains non-zero intensity background noise

# Background correction

Many signal-noise models view the observed log-intensities as a combination of true signal and background noise.



Signal    +    Noise    =    Observed

The 'signal + background' model for the intensities:

$$Y_{ij} = S_{ij} + BG_{ij}$$

$\Rightarrow$ $Y_{ij}$ is the intensity of sample *i* and feature *j*.

$\Rightarrow$ $S_{ij}$ and $BG_{ij}$ are independent random variables.

$\Rightarrow$ $BG_{ij} \sim \mathcal{N}(\mu_i, \sigma_i^2)$

$\Rightarrow$ $S_{ij} \sim \exp(\alpha_i)$

# Background correction

Estimation of $\mu_i$ :

• Fit a density to the $Y_{ij}$, using a kernel density estimator.

• Estimate $\mu_i$ by the mode of the density.

Estimation of $\alpha_i$ and $\sigma_i$ :

$$\hat{\alpha}_i = \sum_{j=1}^p (Y_{ij} - \hat{\mu}_i) I_{\{Y_{ij} \geq \hat{\mu}_i\}} / \sum_{j=1}^p I_{\{Y_{ij} \geq \hat{\mu}_i\}}$$

$$\hat{\sigma}_i^2 = 2 \sum_{j=1}^p (Y_{ij} - \hat{\mu}_i)^2 I_{\{Y_{ij} \geq \hat{\mu}_i\}} / \left( \sum_{j=1}^p I_{\{Y_{ij} \geq \hat{\mu}_i\}} - 1 \right)$$

The background corrected intensity is $B(Y_{ij}) = \mathbb{E}(S_{ij} \,|\, Y_{ij})$ with

$$\mathbb{E}(S_{ij} \,|\, Y_{ij}) = \int_0^\infty s P(Y_{ij} \,|\, S = s) P(S = s) / P(Y_{ij}) ds$$

and

$$f_Y(Y_{ij}) = \int_{-\infty}^\infty f(Y_{ij} - z) f_{BG}(z) dz$$

# Background correction

*Important*

There is no a priori justification for the presented 'signal + background' model (or any of its competitors):

*its usefulness must prove itself in application!*

For instance, check distributional assumptions.

qq-plot

# Normalization

*Motivation for normalization*

Normalization is required to correct for experimental artifacts while preserving the true biological signal.

Normalization balances intensities
→ between dyes, and
→ between hybridizations
in order to allow comparison of gene expression across hybridizations.

- Dyes: within-array normalization.
- Hybridizations: between-array normalization.

# Normalization

Conceptually, normalization adjusts intensities relative to intensities of reference genes whose levels are assumed to be constants between samples.

A set of genes that are to function as reference genes in the normalization must be chosen.

| Genes for normalization | |
|---|---|
| 1. | All genes on the array |
| 2. | Constantly expressed genes |
| 3. | Controls |
| 4. | Rank invariant genes |

# Normalization

*1. All genes on the array*
All genes on the array are used in normalization.

This is sensible when:

a)  only a relatively small proportion of the genes will vary significantly in expression between mRNA-samples, or

b)  there is symmetry in the expression levels of the up/down-regulated genes.

# Normalization

*3. Controls*
Spiked-in controls are synthetic DNA sequences (complementary to oligonucleotides on the array) and included in the mRNA samples at equal amount and should have equal intensities across hybridization.



array 1                    array 2

spiked-in controls

# Normalization

*Within-array normalization* aims

→ to balance intensities of the two dyes, as well as

→ to elimate other systematic differences due to unequal experimental conditions.

Systematic deviations from the line log(G)=log(R) indicate a dye-effect that is to be eliminated by normalization.

# Normalization

*MA plot*
Within-array normalization uses the MA-plot to identify artifacts and detect intensity-dependent patterns in log-ratio's $M_j$.

Statistically, *within-array normalization* subtracts a function $g(\bullet)$ from the individual intensity log-ratio's $M_j$. The function $g(\bullet)$ is computed per array.



before

$M = \log_2 R - \log_2 G$

after

$M = \log_2 R - \log_2 G$

$A = (\log_2 R + \log_2 G)/2$

# Normalization

*Operationalization of between-array normalization*

A transformation of the individual intensities values (or log-ratio's) such that the intensities are comparable across arrays.

→ A transformation is constructed for each hybridization.

→ The functional form of the normalizing transformation is determined by the type of normalization.

*Discuss: scale and quantile normalization*

# Normalization

*Scale normalization*

Assume log-ratios from array *i* follow $\mathcal{N}(0, a_i^2\,\sigma^2)$ .

The scale factors $a_i$ are robustly estimated by:

$$\hat{a}_i \;\; = \;\; \frac{\mathrm{mad}_i}{\sqrt[n]{\prod_{i=1}^{n} \mathrm{mad}_i}}$$

where $\mathrm{mad}_i$ is the median absolute deviation of array *i*:

$$\mathrm{mad} \;\; = \;\; \mathrm{median}_i\{|X_i - \mathrm{median}_j(X_j)|\}$$

Scale normalization is then achieved by dividing the log-ratios by the estimate of $a_i$.

# Normalization

*Quantile normalization* aims to make the distribution of probe intensities the same across arrays. This operationalization is motivated by the assumption that the amount of mRNA in each sample is roughly the same.

It transforms the data from all arrays such that the transformed data follow the $n$-dimensional identity line in the $n$-dimensional qq-plot.

*Rationale*
The quantiles of two identical distribution line up on the diagonal of a qq-plot. This suggests that two datasets could be given the same distribution by equalling their quantiles.

# Normalization

Different distributions ...



… same distribution

# Expression measure

*Operationalization of expression*

An expression measure is a number reflecting the amount of RNA in the sample.

For dual channel arrays:
- the expression measure is simply the log-ratio's $M_j$.

For the Affymetrix single channel array:
- an expression measure is determined by summarizing the probe level data of a probe set (set of features interrogating the same gene) into one number.
  *(not discussed)*

# Expression measure

Expression matrix

Covariate information for samples

|  | Sample 1 | Sample 2 | … | Sample *n* |
|---|---|---|---|---|
| Probe 1 | | | | |
| Probe 2 | | | | |
| ⋮ | Expression signature → | | | |
| Probe *m* | | | | |

Covariate information for probes

Expression profile

# Rubbish?

"Microarrays are the closest thing to fraud we accept in science."

-- *????, ????*

- Inherently noisy.
- Many sources of variation.
- Many preprocessing steps, with lots of arbitrary choices.

# Rubbish?

*Affymetrix spike-in experiment*

→ 14 gene groups are spiked-in at varies concentrations in accordance with a latin square design.

→ Each hybridization has been replicated at least three times.

→ In total 59 hybridization.

→ Array type: HG-U95.

*A proof of principle?*

# Rubbish?



37777_at

36085_at

# Other data have similar issues

*Twitter data*
E.g. can one identify one's political preference on the basis of his/her tweets?

Sometimes easy:


Donald J. Trump
@realDonaldTrump
Follow
Democrats are the problem. They don't care

but meaning not always obvious:


Donald J. Trump
@realDonaldTrump
Follow
Despite the constant negative press covfefe

# Other data have similar issues

*Twitter data*

Harvesting: which tweets to select?

Decide upon:
→ Time period of tweets.
→ Original tweet only?
→ Include retweets?
→ Include replies?
→ Which language?
→ Users' geographical location.
→ Include meta-data like user profile?

# Other data have similar issues

*Twitter data*

Preprocessing issues:

→ URLs, @, #, emoticons, and other symbols,

→ Spelling:

  → 1000, 1,000, 1000.00, 1,000.00, or thousand, or

  → colour or color,

→ Synonyms:

  → loud or noisy (in e.g. a Tripadvisor review),

→ Acronyms:

  → POTUS, LOL, BFF.

→ Remove low frequency words?

→ Remove stop words like "and"?

→ Combine tweets?

An example:
the big promise

# An example

**Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications**

Therese Sørlie[a,b,c], Charles M. Perou[a,d], Robert Tibshirani[e], Turid Aas[f], Stephanie Geisler[g], Hilde Johnsen[b], Trevor Hastie[e], Michael B. Eisen[h], Matt van de Rijn[i], Stefanie S. Jeffrey[j], Thor Thorsen[k], Hanne Quist[l], John C. Matese[c], Patrick O. Brown[m], David Botstein[c], Per Eystein Lønning[g], and Anne-Lise Børresen-Dale[b,n]

Using 78 breast cancer profiles, five subtypes are distinguished:
- Basal
- ERBB2
- Luminal A
- Luminal B
- Normal

# An example

*Traditional medicine*

All treated with the same drug

Standard treatment may not be beneficial to everyone.

Subgrouping of breast cancers suggest patients from different subgroups may need different therapy.

# An example

*Personalized medicine*
Individualized treatment based on patient's genetic characteristics.

# An example

Why do people believe these breast cancer subtypes?
1) Subtypes exhibit different clinical outcome.

# An example

Why do people believe these breast cancer subtypes?
2) Exhibit different morphology.

# An example

Why do people believe these breast cancer subtypes?
3) Subtypes have been confirmed.

## Repeated observation of breast tumor subtypes in independent gene expression data sets
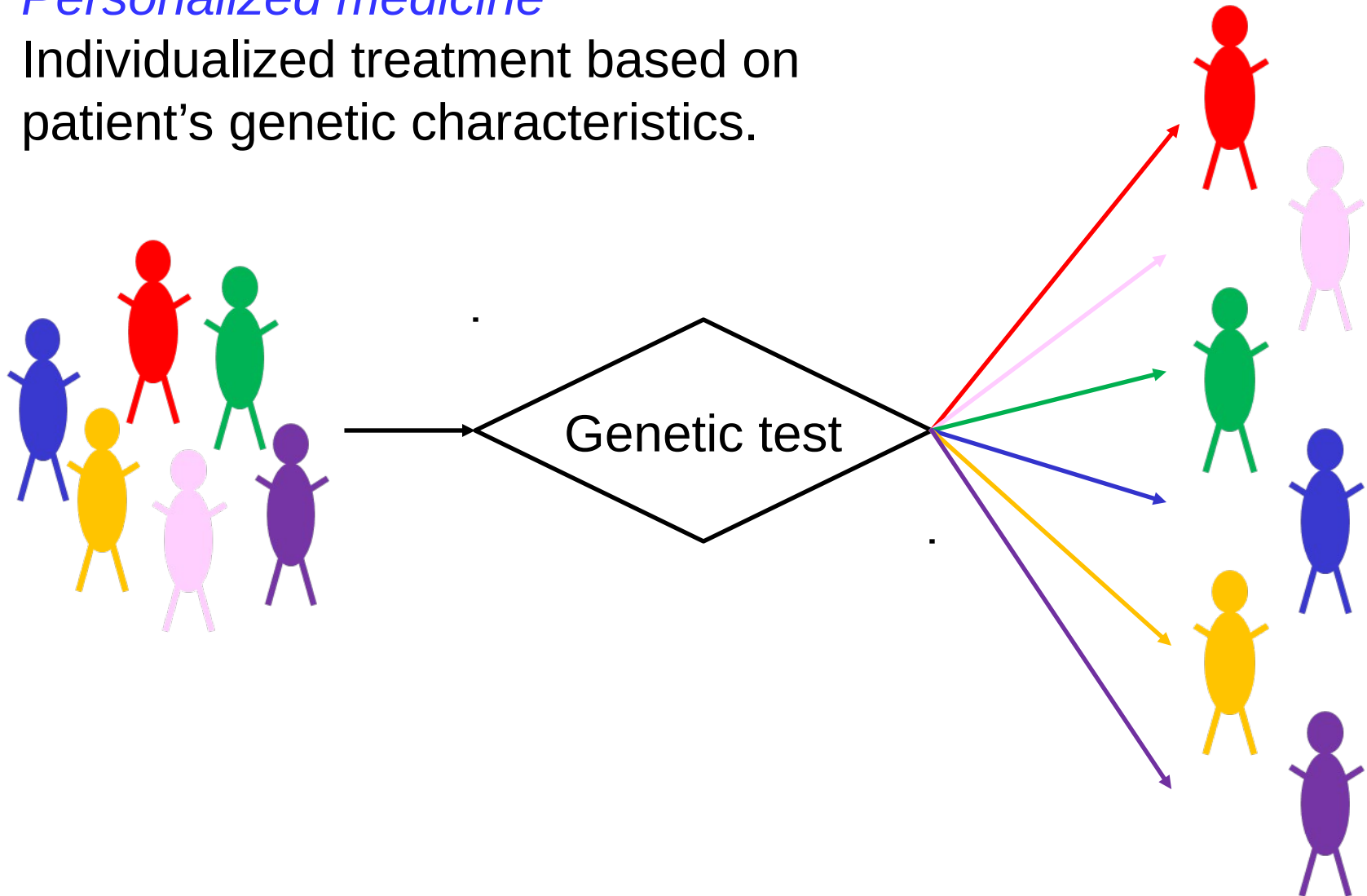
Therese Sørlie*, Robert Tibshirani[†], Joel Parker[‡], Trevor Hastie[§], J. S. Marron[¶], Andrew Nobel[¶], Shibing Deng[‖], Hilde Johnsen**, Robert Pesich*, Stephanie Geisler[††], Janos Demeter*, Charles M. Perou[‡,‡‡], Per E. Lønning[††], Patrick O. Brown[§§], Anne-Lise Børresen-Dale**, and David Botstein*[¶¶]

# An example

Medio 2012, the story continues …

## ARTICLE

## The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups

Christina Curtis[1,2]†*, Sohrab P. Shah[3,4]*, Suet-Feung Chin[1,2]*, Gulisa Turashv:
Doug Speed[2,5]†, Andy G. Lynch[1,2], Shamith Samarajiwa[1,2], Yinyin Yuan[1,2], Ste
Ali Bashashati[3], Roslin Russell[2], Steven McKinney[3,4], METABRIC Group‡, Anit
Gordon Wishart[8], Sarah Pinder[9], Peter Watson[3,4,10], Florian Markowetz[1,2], Lei
Anne-Lise Børresen-Dale[6,12], James D. Brenton[2,13], Simon Tavaré[1,2,5,14], Carlo:

Inclusion of more molecular information suggests the existence of 10 subgroups.

# An example

How many subgroups really exist?

Genetically, everybody is unique. Thus …



19??  —  2002  —  2012  —  20??

8 billion?
…
personalized medicine
to the max?

# The curse of dimensionality

# Curse of dimensionality

*i)* The high-dimensional space is enormous and data points are isolated.

*Unit cubes*

p=1

p=2

p=3

*Maximum distances*

$$\sqrt{1^2} = 1 \qquad \sqrt{1^2 + 1^2} = \sqrt{2} \qquad \sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}$$

# Curse of dimensionality

Histograms of Euclidean distances between random vectors from the *p*-dimensional unit cube.



Distances grow with *p* and (relatively) more homogeneous.

*Question*
What does this do to the "nearest neighbor" concept?

# Curse of dimensionality

*ii)* Volume of unit balls vanishes as *p* increases.

The volume of a sphere with radius *r* in a *p*-dimensions:

$$V_p(r) = [\Gamma(p/2 + 1)]^{-1} \pi^{p/2} r^p$$

For $p = 20$, $V_p(1) = 1.73 \times 10^{-13}$.



*r*=1

*Question*
What are the consequences for the sample size?

# Curse of dimensionality

A well-designed experiment provides good coverage of the design space.

→ Design space: unit cube $[0,1]^p$.
→ Distribute $n$ points s.t. union of unit balls around the points encompass the unit cube.
→ Calculate required $n$ for varying dimensions $p$.

| $p$ | $n$ |
|---|---|
| 20 | 39 |
| 30 | 43630 |
| 50 | $5.7 \times 10^{12}$ |
| 100 | $4.2 \times 10^{39}$ |
| 150 | $1.28 \times 10^{72}$ |

# Curse of dimensionality

*iii)* Volume of unit balls concentrates around the crust.

Denote *p*-dimensional ball with radius *r* by: $B_p(0, r)$.

Define the "crust" by:

$$C_p(r) = B_p(0, r) \setminus B_p(0, 0.99\,r).$$

**ball**       **crust**

Then, the crust-to-ball volume is:

$$\text{Vol}[C_p(r)] \,/\, \text{Vol}[B_p(0, r)]$$

# Curse of dimensionality

The univariate normal distribution concentrates most mass around its mean and has thin tails.

This is reversed for large *p*:

$$P[\exp(-\tfrac{1}{2}\|\mathbf{X}\|_2^2) \geq \delta] \leq (\delta 2^{p/2})^{-1}$$



*Consequence*
Rare events may not be so rare.

This lecture series

# Issues: multiplicity

A comparison of the expression levels of gene *A* between two groups:

```
        Welch Two Sample t-test

data:  group 1 and group 2
t = -8.6449, df = 17.284, p-value = 1.099e-07
alternative hypothesis: true difference in means is not equal to 0
```

- This is a rather small *p*-value.

- Getting such a small *p*-value is unlikely.

- Is it still unlikely if we acknowledge that the gene is one of 40000 on the microarray?

# Issues: multiplicity

*The multiplicity problem*
Each individual test has a specified type I error probability. This probability of committing a type I error increases with the number of tests.

The probability of at least one false positive finding in m tests is given by:

$$1 - (1 - P(\sim H_0 \mid H_0))^m$$

$$=$$

$$1 - (1 - \alpha)^m$$

| m | $1 - (1 - P(\sim H_0 \mid H_0))^m$ |
|---|---|
| 1 | 0.0500 |
| 2 | 0.0975 |
| 5 | 0.2262 |
| 10 | 0.4013 |
| 100 | 0.9941 |

# Issues: multiplicity

Decreasing the rejection level reduces the probability of a false positive.

# Issues: multiplicity

*Problem*
→ many traits, many tests,
→ large number of false positives.

*Multiple testing*
→ generalization of type I error,
→ control of this generalized type I error,
→ control of number of false positives.

*Techniques (in lecture series)*
→ FWER,
→ FDR.

# Issues: shrinkage

Estimation of the variance of expression levels of gene *A*, with only few samples available.



*sample size vs. confidence*

Few samples → large uncertainty.

# Issues: shrinkage

Additional information available:
variance estimates of 40000 other genes.



*histogram of variance estimates*

Confidence interval of overall pooled variance estimate:
very, very small.

# Issues: shrinkage

Individual variance estimate:
→ unbiased, but large uncertainly

Overall pooled variance estimate:
→ biased, but very low uncertainty

*Why not exploit the strengths of boths?*
E.g. by combined estimator:

$$(1-\theta)\ s^2_{individual} + \theta\ s^2_{overall}$$

The individual estimator is "shrunken" towards the overall.

# Issues: shrinkage

*Problem*
→ low sample: highly variable estimates,
→ low-reproducibility


*Shrinkage*
→ traits are "comparable",
→ borrow information across traits,
→ stabilizes estimation and improve inference.


*Techniques (in lecture series)*
→ Stein estimator,
→ Empirical Bayes.

# Issues: penalized estimation

*A common objective*
predict clinical outcome from gene expression levels.

Data available:
→  a few hundred samples at best,
→  # covariates ≈ 40000.

Harrell (2001) gives the following rule-of-thumb:
> *For each continuous covariate in the model 10–20*
> *observations are needed to detect reasonably sized*
> *effects with reasonable power.*
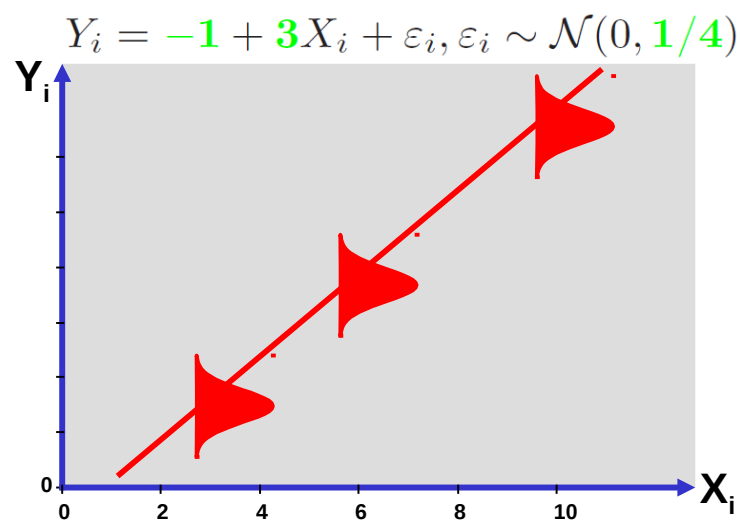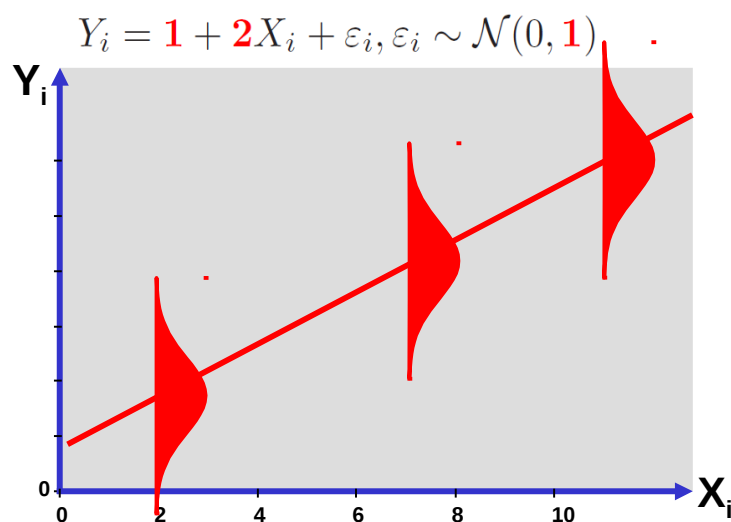
Where does this put us? A model with 20 genes?

# Issues: penalized estimation

*Identifiability*
A statistical model is *identifiable* if for any two choices of the parameter $\theta_1$ and $\theta_2$, such that $\theta_1 \neq \theta_2$, the resulting probability distributions differ: $P_{\theta_1} \neq P_{\theta_2}$.

*Fact*
The linear regression is identifiable.

$Y_i = \mathbf{1} + \mathbf{2}X_i + \varepsilon_i, \varepsilon_i \sim \mathcal{N}(0, \mathbf{1})$

$Y_i = \mathbf{-1} + \mathbf{3}X_i + \varepsilon_i, \varepsilon_i \sim \mathcal{N}(0, \mathbf{1/4})$

# Issues: penalized estimation

*Identifiability*

Empirically, when p > n, the parameters cannot uniquely be identified from the data. That is, multiple parameter choices yield the same model.

Data available:
→  100 samples,
→  10000 covariates,

Then:
→ Let the first 100 covariates be linearly independent.
→ The same holds for the second 100 covariates.
→ Both sets of covariates produce a linear regression model with a perfect fit.

How do we distinguish between the two?
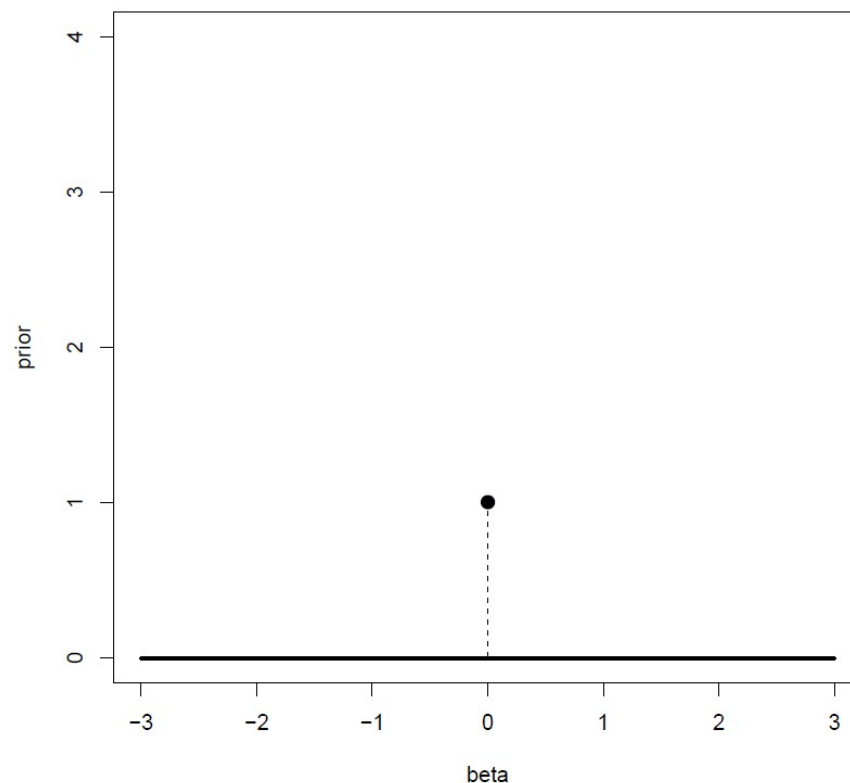
# Issues: penalized estimation

Additional information may help.

*In extremis*
Would one know which (max) 20 covariates to include,
Harrell (by his rule-of-thumb) would not object.

A natural way to include
such information is e.g.
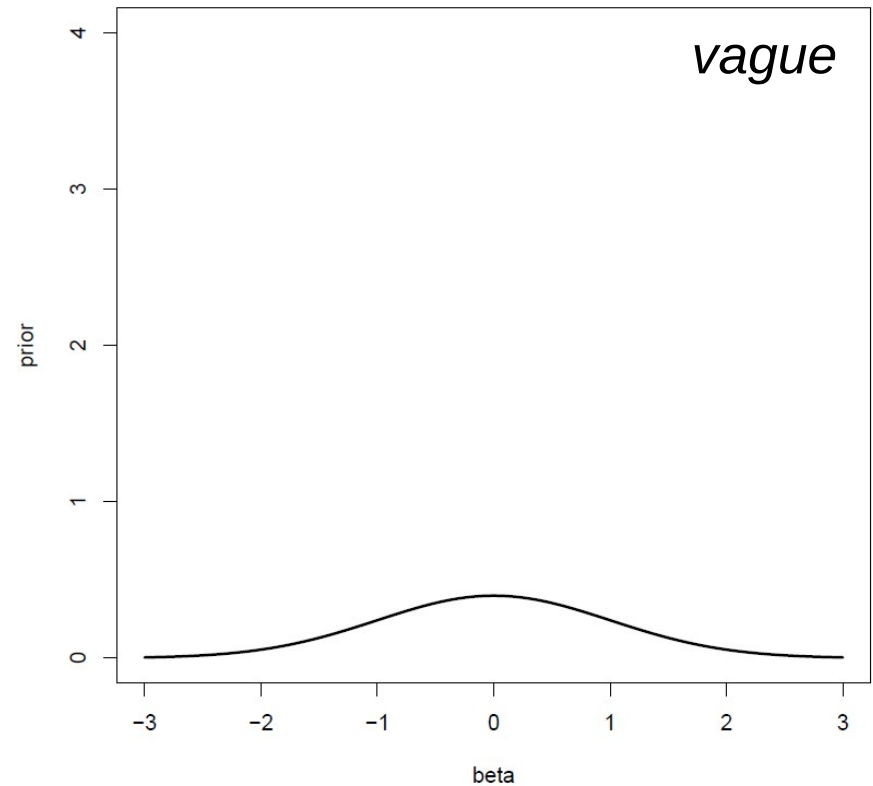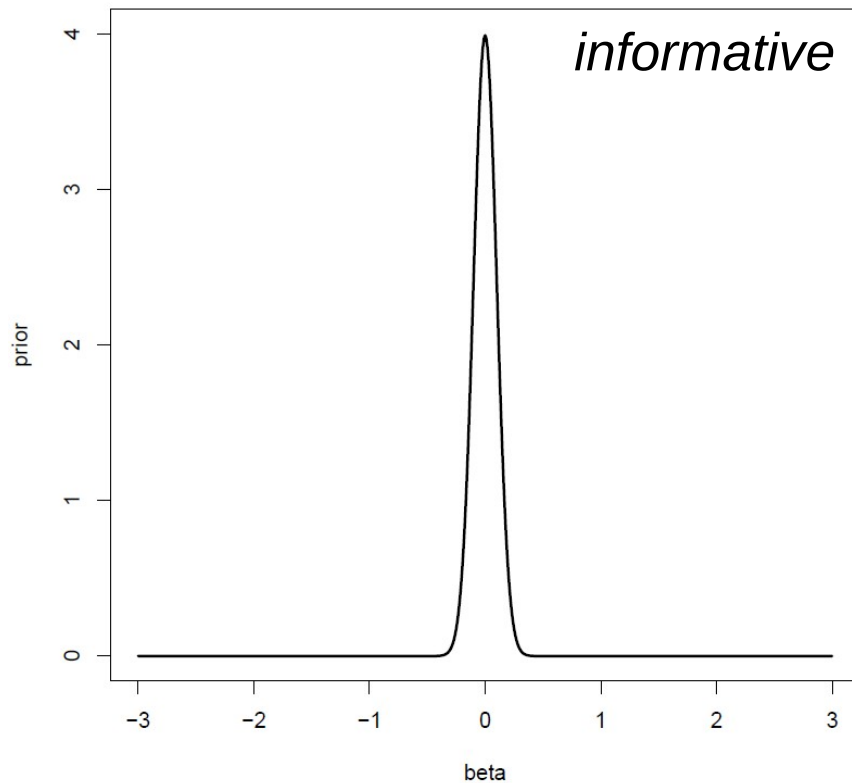through the specification of
a prior.

A *prior with a point mass* at
zero omits the covariate
from the model.

# Issues: penalized estimation

Often no knowledge on relevant covariates. Data may help in the selection of the prior. E.g., very large sample size: no informative prior needed.

# Issues: penalized estimation
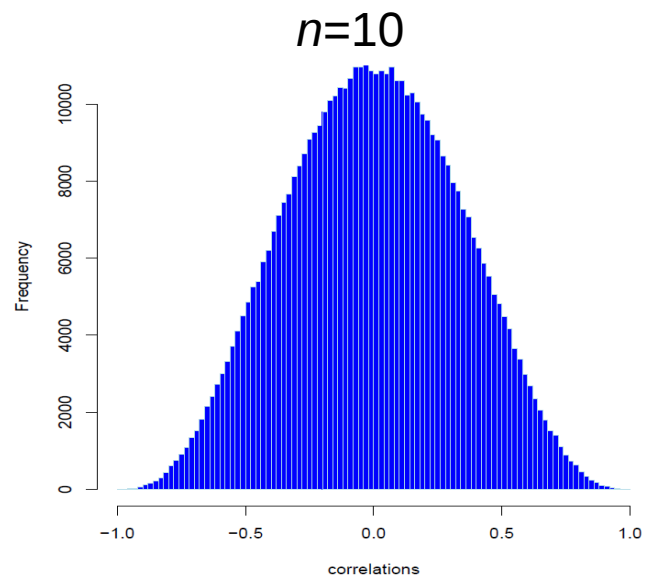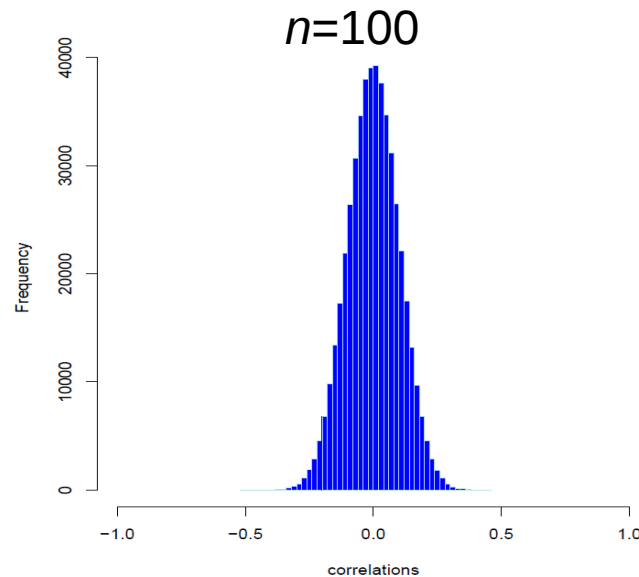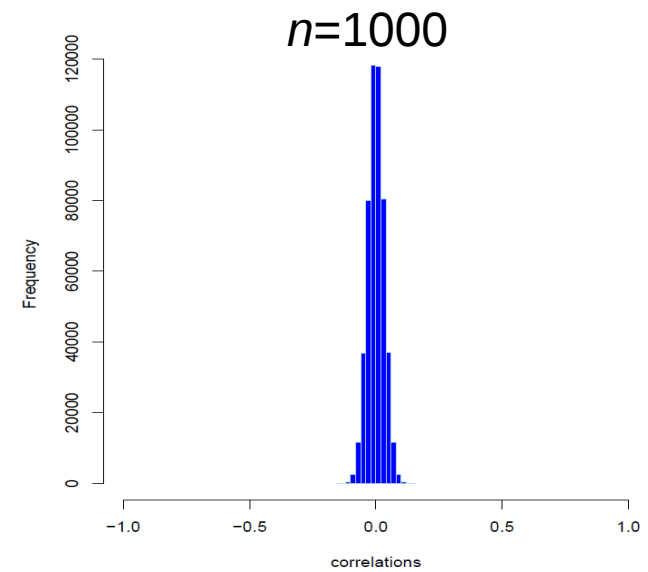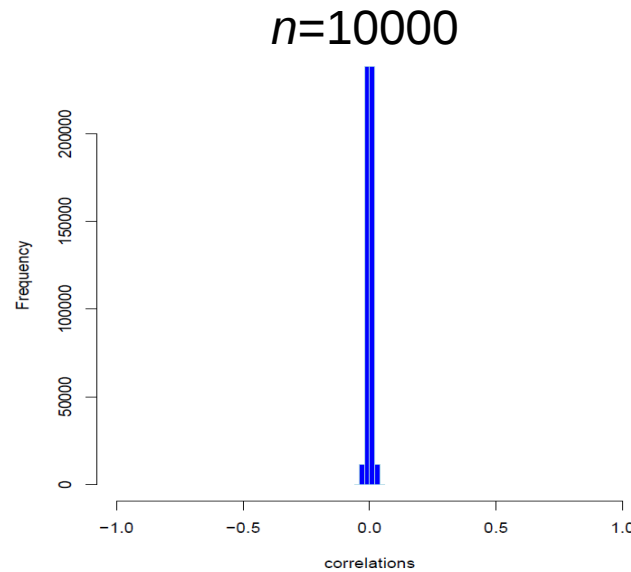
*Problem:*
Correlation
estimates
inflate

Data:
→ $Y_i \sim N(0_p, I_{pp})$
→ $p = 1000$
→ $n$ varies

Calculate
sample
correlation
matrix.

# Issues: penalized estimation

*Problem*
→ few samples, many parameters,
→ estimation is frustated by $p > n$.

*Penalized estimation*
→ estimation procedures for $p > n$,
→ less variable but biased estimates,
→ illustrated on reconstruction of networks.

*Techniques (in lecture series)*
→ ridge regression,
→ lasso regression,
→ ridge and lasso (inverse) covariance estimation.

# Issues: asymptotics (not treated)

*Problem:*

→ when *p* grows, parameters no longer fixed, e.g.:

$$\mathcal{N}_p(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$$

*Various asymptotic limits discerned:*

→ $n \succ p$: $n \to \infty$ and $p$ fixed ,

→ $n \succeq p$: $n \to \infty$ and $d = \mathcal{O}(n)$,

→ $p \succeq n$: $p \to \infty$ and $n = \mathcal{O}(p)$,

→ $p \succ n$: $p \to \infty$ and $n$ fixed .

# Big data vs. high-dimensional data

# Big vs. high

Big data:
→ large $n$: many individuals, large sample size.
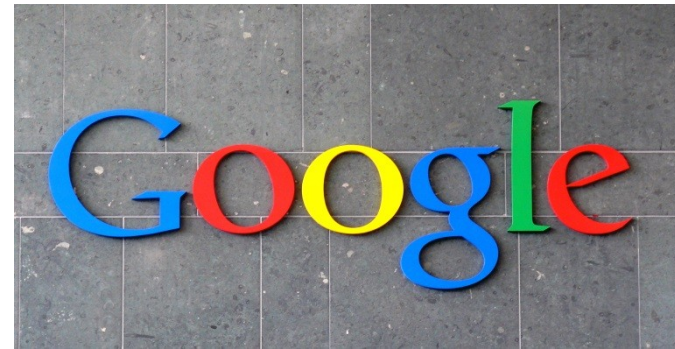→ large $p$: information on many traits of these individuals.

Google:
→ $n$ large: many people use Google software.
→ $p$ large: Google registers everything these $n$ do.
Similarly, Facebook, ING, et cetera.

Why:
→ many individuals available.
→ information is cheap to gain.

# Big vs. high

High-dimensional data:
→ small $n$: few individuals, small sample size.
→ large $p$: information on many traits of these individuals.

VUmc:
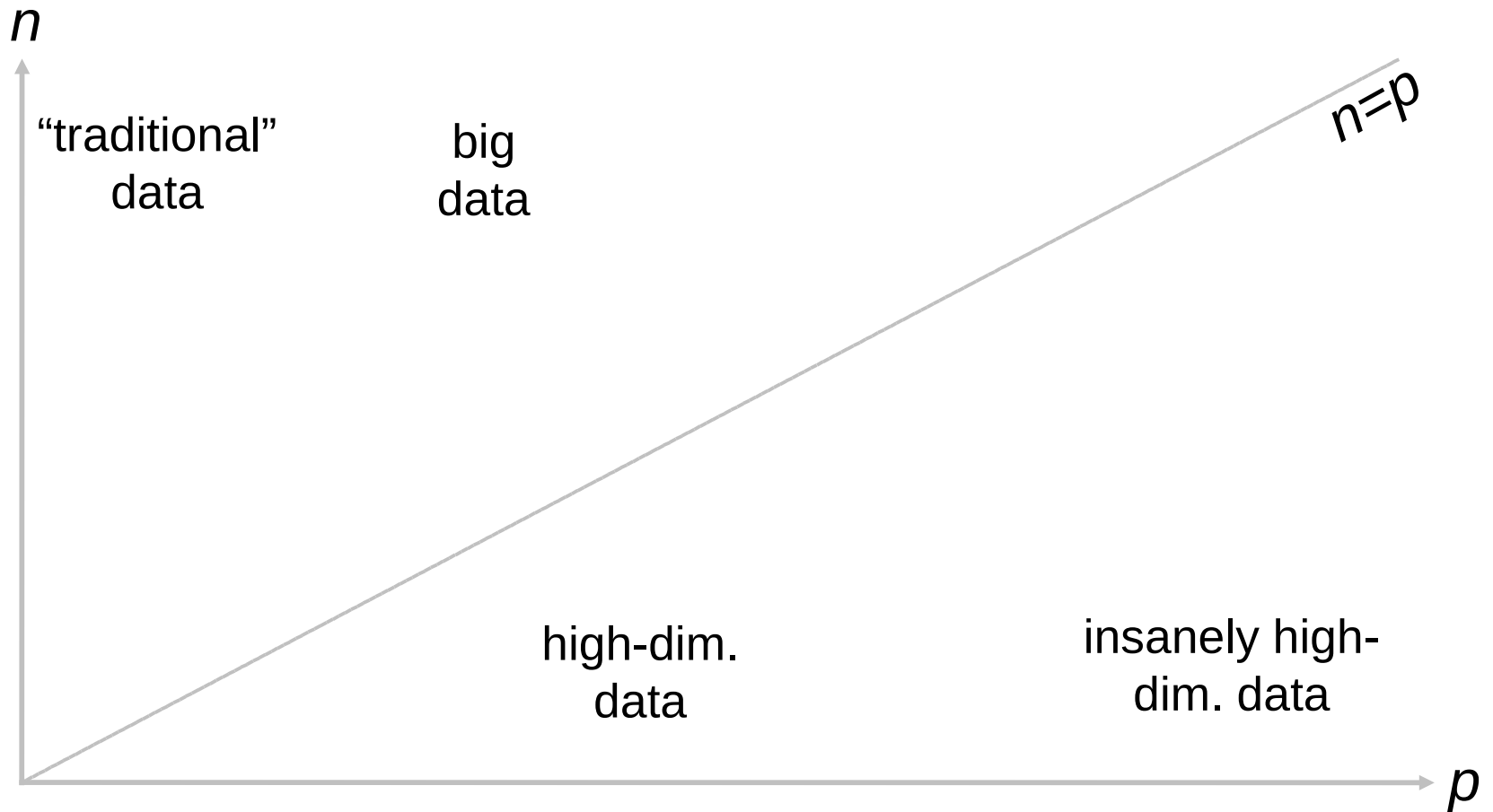→ $n$ small: few cancer patients.
→ $p$ large: many traits.

Why:
→ individuals with particular disease not abound.
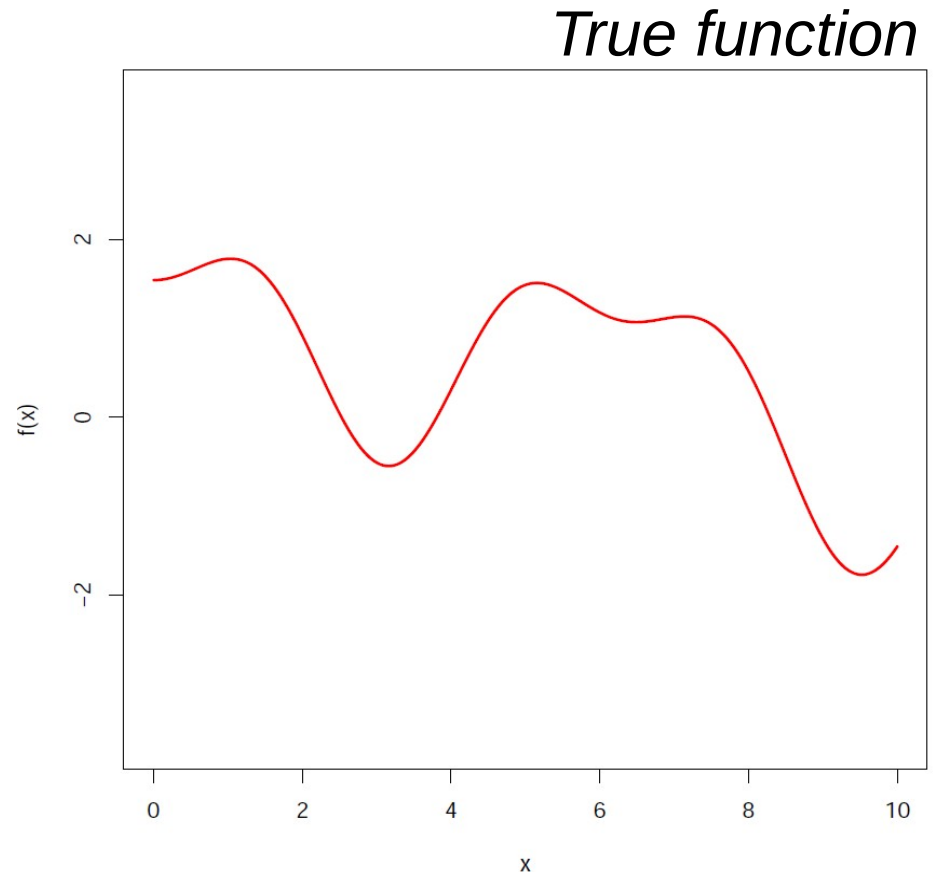→ information is expensive.

**VUmc**

# Big vs. high

# Big vs. high

Nonlinear function:
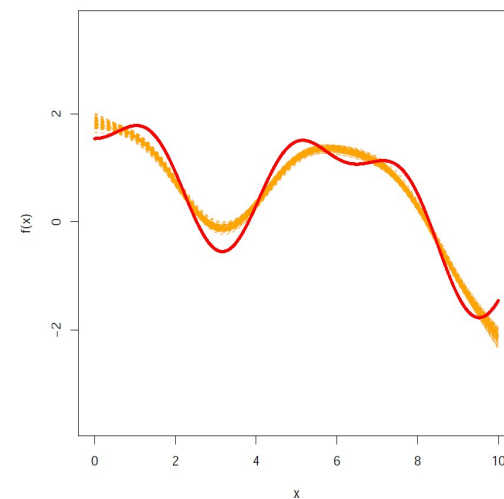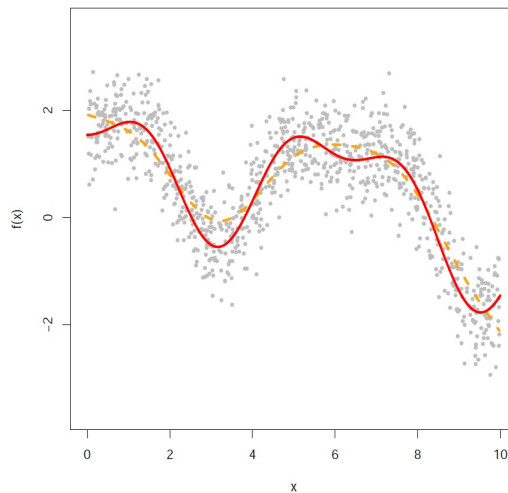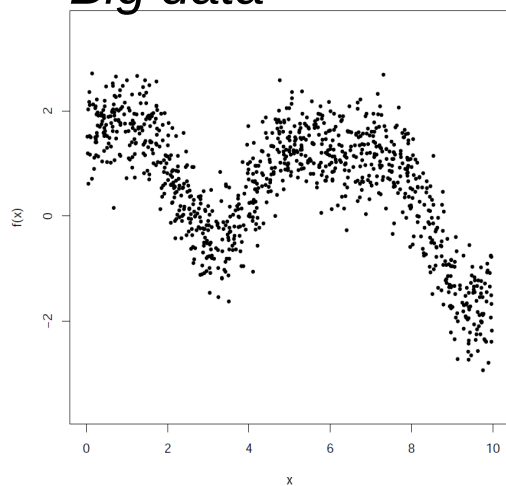$f(x) = \sin(x) + \sin(x^2) + \cos(\exp(x^2/100))$.

Estimate f(x) from:
- → big data ($n$=1000),
- → high-dim. data ($n$=10).

Approximate f(x) by
spline with $p$ degrees of
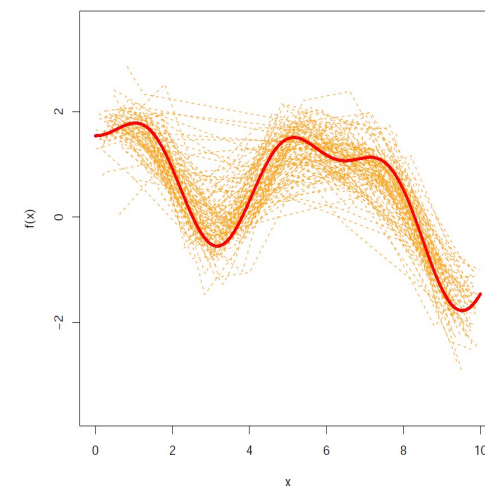freedom. Use $p < 7$, while
f(x) requires $p \gg 7$.

*True function*

# Big vs. high

*Big data*



*High-dimensional data*

# Big vs. high

Big and high-dimensional data often differ in the experimental design underlying the data.

**Big data**

*Practice*
Google collects virtually anything it can gets its hands on.

*Design*
Observational at best

**High-dim. data**

*Practice*
VUmc financial sources limited: careful planning of experiments.

*Design*
Observational, but often well-controlled experiments.

# Big vs. high

Big and high-dimensional data often used for different purposes.

Big data

*Practice*
Google optimizes advertisement revenue.

*Aim*
<u>Predict</u> behaviour of the internet user.

High-dim. data

*Practice*
VUmc tries to cure patients.

*Aim*
<u>Understanding</u> of the disease mechanism.

# Big vs. high

Both big and high-dimensional data sometimes originate from diffuse sources.

**Big data**

Google measures all, but may also acquire third-party data.

**High-dim. data**

VUmc measures molecular and clinical traits of patients.

The data thus comes from various sources, possibly in different formats and with varying quality. To benefit from this multitude of sources is challenging.

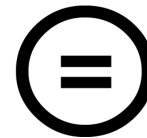# References

# References

Bolstad, B.M., Irizarry, R.A., Astrand, M., Speed, T.D. (2003), "A comparison of normalization methods for high density oligonucleotide array data based on variance and bias", *Bioinformatics*, **19**(2), 185-193 .

Cleveland, W.S. (1979), "Robust locally weighted regression and smoothing scatterplots", JASA, 74(368), 829-836.

Giraud, C. (2015), *Introduction to High-Dimensional Statistics*. CRC Press.

Harrell, F.E. Jr. (2001), *Regression Modelling Strategies, with Applications to Linear Models, Logistic Regression, and Survival Analysis*. New York: Springer.

Irizarry, R.A., Hobbs, B., Collin, F., Beazer-Barclay, Y.D., Antonellis, K., Scherf, U., Speed, T.D. (2003), "Exploration, normalization, and summaries of high density oligonucleotide array probe level data", *Biostatistics*, **4**(2), 249-264.

Koch, I. (2013), *Analysis of Multivariate and High-dimensional Data,* Cambridge University Press.

Mayer-Schönberger, V., Cukier, K. (2013), *Big Data: A Revolution That Will Transform How We Live, Work and Think*, John Murray Publishers.

Nguyen, D.V., Bulak Arpat, A., Wang, N., Carroll, R.J. (2002), "DNA microarray experiments: biological and technological aspects", *Biometrics*, 58, 701-717.

Quackenbush, J. (2002), "Microarray data normalization and transformation", *Nature Genetics Supplement*, **32**, 496-501.